



Molecular Modeling

Prediction of Protein 3D Structure from Sequence

Vimalkumar Velayudhan

Jain Institute of Vocational and Advanced Studies

May 21, 2007



1 Introduction

- What is Molecular Modeling?
- Methods in Molecular Modeling

2 Homology Modeling

- Steps
- Guidelines
- Modeling Programs

3 Further Reading



What is Molecular Modeling?



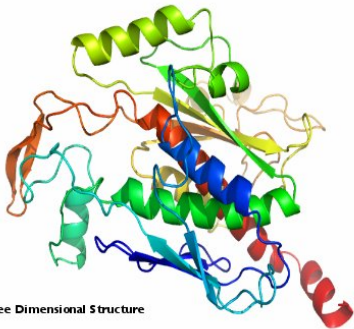
- The Prediction of a protein's three dimensional structure from its sequence
- A computational method based on our understanding of protein structures

Principle

The Sequence of a protein ie., the order of amino acids determine the 3D structure of the protein and hence, its function

Amino acid
sequence

Gly Phe Leu Gly Ala Ala Gly Ser Thr Met Gly Ala



Three Dimensional Structure



Methods in Molecular Modeling

- 1 Homology Modeling¹
- 2 Threading of Fold Recognition
- 3 *ab initio* Prediction

Accuracy

The accuracy of the methods are in the following order
Homology Modeling > Threading > *ab initio* Prediction

¹Also known as Comparative Modeling

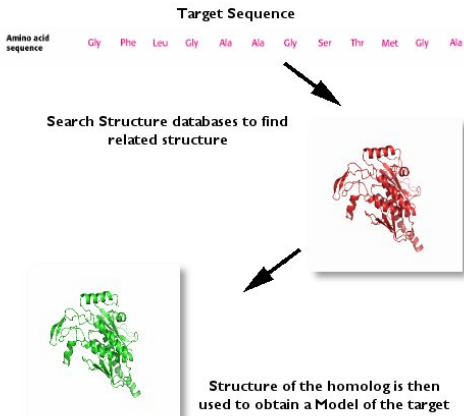


- Predicting the structure of a protein using the structure of its homolog²
- The protein whose structure is not known is referred to as the target and the structure of the homolog is referred to as the template

Basic Principle

Similar Sequences tend to adopt Similar Structures

²Proteins which share a common ancestor in evolution





Given the sequence of a protein, its structure can be predicted in the following steps

- 1 Template detection
- 2 Target –Template Alignment
- 3 Backbone generation
- 4 Modeling of Side-chains and Loops
- 5 Model Validation and Optimization



Template detection



- A template is a protein whose structure is already known³
- Target refers to the protein we would like to model - whose structure is unknown

Template Detection refers to the identification of a suitable template corresponding to the target by database similarity searches

- Can be performed by doing a BLAST⁴ search⁵ against the PDB⁶ database

³ by X-ray Crystallography, NMR or other techniques

⁴ A program to search a database with a sequence to identify related sequences

⁵ blastp - Protein-Protein BLAST

⁶ Protein Data Bank - a worldwide repository for protein structures



Template detection



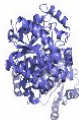
Target Protein Sequence (Obtained from database like Swiss-Prot)

```

|CYB_HUMAN P00156 Cytochrome b.
MTPMRKINPLMKLINHSFIDLPTPSNISAWNFGSLLGACILQITTGLFLAMHYS PDAS
TAFSSIAHITRDVNYGWIIRYLHANGASMFICFLHIGRGLYYGSFLYSETWNIGIILL
LATMATAFMGYVLPWGQMSFWGATVITNLLSAIPYIGTDLVQWIWGGYSVDSPTLTRFFT
FHFILPFIIAALATLHLLFLHETGSNNPLGITSHSDKITFHPYYTIKDALGLLLFLLSLM
TLTLFSPDLLGDPDNYTLANPLNTPPHIKPEWYFLFAYTILRSVPNKLGGVLALLLSILI
LAMIPILHMSKQOSMMFRPLSQSLYWLLAADLLILTWIGGQPVSYPTIIGQVASVLYFT
TILILMPTISLIENKMLKWA
  
```



**Search PDB to find
related Structures
(potential templates)**





Template detection



- Atleast 25% sequence identity is required between the query and the subjects from PDB
- A number of other parameters are also involved in the selecting the right template
 - Resolution - the higher⁷ the better
 - The template should cover the entire length of the target
 - Gaps should be minimal

⁷Higher resolution will correspond to a low numeric value - Ex., 1.0 is better than 2.0



Target-Template Alignment



- An alignment is necessary to state which residues in the target correspond to which residue in the template
- The model-building program uses this information to build the backbone of the target
- Can be performed using programs like ClustalX, T-Coffee or the Modeling program itself



Target-Template Alignment



```

1      10      20      30      40      50      60
lumOA  M...NTTGRFLRRTTFGESHGIVICGVLDGMPSGIRIDYALLNEMKRRRGGRNVFTPRRDKVITSG
SaLPA  MAGNTTGRFLRRTTFGESHGIALICLVDGVPPIGITEADLQNDLDRRPGESRYVTCRRRPFDDVITLSC

70     80     90     100    110    120    130
lumOA  VFEDFSTGTPIGLIINCRARSQDYDNINLFRPSHADTYFHKYGRDRFRGGGRSSAREAIRVAAGAF
SaLPA  VFDGVLTGTSLGLIENIDQRSQDYSAIKDVFPRGHADNTYEQKYGRDRFRGGGRSSAREAIRVAAGAF

140    150    160    170    180    190    200
lumOA  AKMLLRRLIGITVCESGIETICGIIKAKNYDFNHALKSEIFALDEEQEEAKTATQNAERNNHDSITGVALLRA
SaLPA  AKKYLRLRKFGLEIRGCLIQMGDIPLEIKDWQVELNPFIFCPDADKLDALDEIMRAARKKEGDSIGAKVTVV

210    220    230    240    250    260    270
lumOA  RSIKTNGKLPVGLGCLYARLDAIYAAAMGILNVKAVEIGKQVESSILKGSRYNDLMDQVGLSNRSGG
SaLPA  AS....GVVAGLGEFVFDRLDAIYAAAMSINAVKVEIGKGFNVVALRGSQNRDLEITACGFQSNHAGG

280    290    300    310    320    330    340
lumOA  VLGGMSNGGEEIVRVHFKPTFSIFCPORTIDINCNECECLLKGRHDPCHAIRGSVVCESLALVLDMMVL
SaLPA  ILGGISGGGHIYAHMALKPTFSITVPGRTINRMGEVLEMITKGRHDPGVGIRAMPYLAALALVLDHML

350    360
lumOA  LNLTSKITEYLKTIYNEN
SaLPA  RHRACNADVKTIPRW.

```

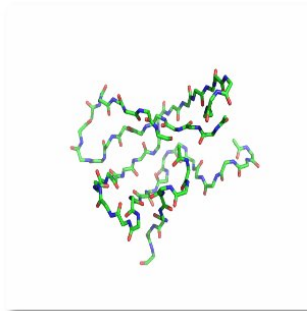


- Based on the information in the alignment, the model-building program generates the backbone of the target
- Backbone refers to the repeating N - C $_{\alpha}$ - C atoms of the polypeptide chain

1 10 20 30 40 50 60
 luxA M...NFGHFRITTTFGESAGDVTGGLVLOGNSCTKISQVLLDNEMSRGGGRNVSPRHRDNDVETS
 SalPA MAGNTGGLFRITTTFGESGLA GCTSDGPPGIPITTAQIHDLPARRGGTSRYTTCRRDPRVSL



Using the Alignment,
A backbone of the target
is generated





Modeling of Side-Chains and Loops

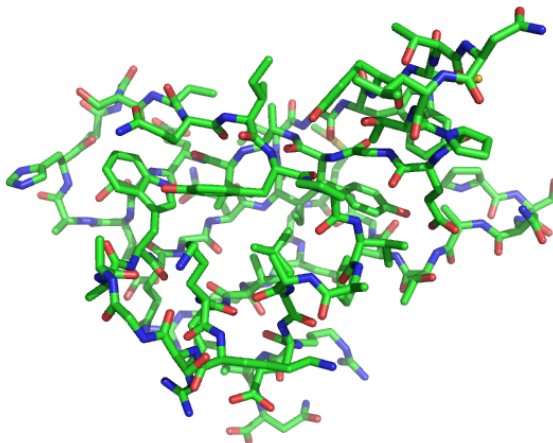


- Once the backbone is built, the side-chains of the residues are then modeled by
 - Using the information from the alignment if the alignment is conserved
 - If residues are different, an inbuilt library of side-chain conformations⁸ is consulted
- The same principle is applied to model loop regions too

⁸Referred to as Side-chain rotamer library



Modeling of Side-Chains and Loops





Model Validation and Optimization



- After the backbone, sidechains and loop regions have been modelled
- The model is checked for its quality
 - Done by assessing parameters like bond-lengths, bond-angles etc.,
- If there are problems at particular regions, the alignment can be verified and adjusted if required
- The steps are performed once again until a satisfactory model is obtained

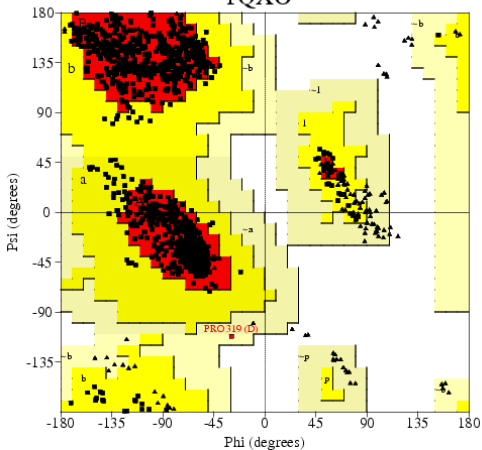


Model Validation and Optimization



PROCHECK

Ramachandran Plot 1QXO



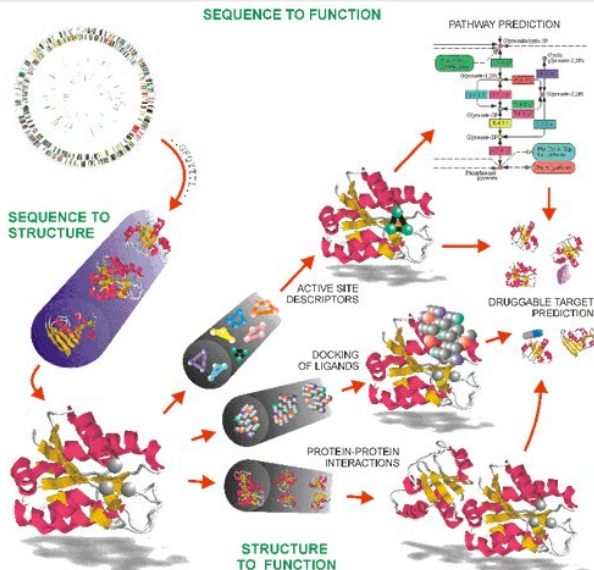


Alignment The most important step! - incorrect alignment will result in a model with errors

Quality of the template - good resolution and stereochemical quality

Percentage Identity between the target and template - the higher the better ⁹

⁹25% is the minimum





Modeling Programs



■ Academic

- MODELLER - <http://www.salilab.org>
- Deepview and SWISS-MODEL - <http://www.expasy.org/spdbv>
- CPH-Models - <http://www.cbs.dtu.dk/services>

■ Commercial

- Accelrys Insight II - <http://www.accelrys.com>



- 1 Structural Bioinformatics - Philip E Bourne, John Wiley & Sons Publications
- 2 Bioinformatics - From Genomes to Drugs, Thomas Lengauer, Wiley publications